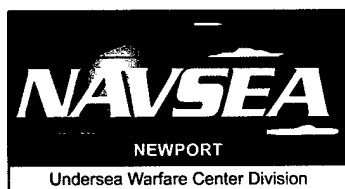


NUWC-NPT Technical Report 11,469
6 November 2003

Joint Probability Density Function of Selected Order Statistics and the Sum of the Remainder as Applied to Arbitrary Independent Random Variables

Albert H. Nuttall
Surface Undersea Warfare Department



20040114 103

**Naval Undersea Warfare Center Division
Newport, Rhode Island**

Approved for public release; distribution is unlimited.

PREFACE

The work described in this report was funded under Project No. A101504, "Automatic Signal Classification," principal investigator Stephen G. Greineder (Code 2121). The sponsoring activity is the Office of Naval Research, program manager John Tague (ONR 321US).

The technical reviewer for this report was Paul M. Baggenstoss (Code 2121).

Reviewed and Approved: 6 November 2003



Donald A. Aker
Head, Surface Undersea Warfare Department



REPORT DOCUMENTATION PAGE

Form Approved
OMB No. 0704-0188

Public reporting for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.

1. AGENCY USE ONLY (Leave blank)

2. REPORT DATE
6 November 2003

3. REPORT TYPE AND DATES COVERED

4. TITLE AND SUBTITLE

Joint Probability Density Function of Selected Order Statistics and the Sum of the Remainder as Applied to Arbitrary Independent Random Variables

5. FUNDING NUMBERS

6. AUTHOR(S)

Albert H. Nuttall

7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)

Naval Undersea Warfare Center Division
1176 Howell Street
Newport, RI 02841-1708

8. PERFORMING ORGANIZATION
REPORT NUMBER

TR 11,469

9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)

Office of Naval Research
Ballston Centre Tower One
800 North Quincy Street
Arlington, VA 22217-5660

10. SPONSORING/MONITORING
AGENCY REPORT NUMBER

11. SUPPLEMENTARY NOTES

12a. DISTRIBUTION/AVAILABILITY STATEMENT

Approved for public release; distribution is unlimited.

12b. DISTRIBUTION CODE

13. ABSTRACT (Maximum 200 words)

A set of N independent random variables $\{x_n\}$ with arbitrary probability density functions $\{p_n(x)\}$ are ordered into a new set of dependent random variables, each with a different probability density function. From this latter set, the $M - 1$ largest random variables are selected. Then, the sum of the remaining $N + 1 - M$ random variables is computed, giving a total of M dependent random variables.

Several statistics are computed for these M random variables, including their joint probability density function, a quantity called the "combined probability and joint probability density function" of particular selections, and their distribution. Most of the results can be expressed as a single Bromwich contour integral in the moment-generating domain. This integral is most easily numerically evaluated by locating (approximately) the real saddlepoint of the integrand and passing the contour through that point. Very high accuracy in the joint probability density function evaluations is available by using numerical integration on this latter contour, instead of a saddlepoint approximation.

14. SUBJECT TERMS

Contour integral
Saddlepoint
Order Statistics

Moment-Generating Function
Joint Probability Density Function
Signal Detection

15. NUMBER OF PAGES
22

16. PRICE CODE

17. SECURITY CLASSIFICATION
OF REPORT
Unclassified

18. SECURITY CLASSIFICATION
OF THIS PAGE
Unclassified

19. SECURITY CLASSIFICATION
OF ABSTRACT
Unclassified

20. LIMITATION OF ABSTRACT
SAR

TABLE OF CONTENTS

	Page
LIST OF ABBREVIATIONS, ACRONYMS, AND SYMBOLS	ii
INTRODUCTION	1
DERIVATIONS OF JOINT PROBABILITY DENSITY FUNCTIONS	2
Largest Random Variable.....	2
Two Largest Random Variables.....	4
Three Largest Random Variables.....	8
Four Largest Random Variables	10
$M-1$ Largest Random Variables	12
Recursion for m -th Order Sum.....	12
Second-Largest Random Variable	14
SUMMARY	16
REFERENCES	16

LIST OF ABBREVIATIONS, ACRONYMS, AND SYMBOLS

a_n	n -th coefficient in exponential densities, equation (5)
A_2	Fundamental two-dimensional sum, equation (53)
C	Contour of integration, equation (9)
CDF	Cumulative distribution function
$c_n(x)$	Cumulative distribution function of \mathbf{x}_n , equation (1)
$c_n(z, \lambda)$	Auxiliary function, equation (3)
EDF	Exceedance distribution function
$e_n(x)$	Exceedance distribution function of \mathbf{x}_n , equation (18)
$e_n(z, \lambda)$	Auxiliary function, equation (60)
M	Number of random variables in set under consideration
MGF	Moment-generating function
n	Index of random variables
N	Number of random variables, equation (1)
PDF	Probability density function
$p_n(x)$	Probability density function of \mathbf{x}_n , equation (1)
$P(z, \lambda)$	Product function, equation (3)
$q_2(m)$	Combined probability and joint probability density, equation (11)
$q_2(z_1, z_2)$	Joint probability density function, equations (15) and (16)
$Q(m)$	Probability that \mathbf{x}_m is the largest random variable, equation (2)
RV	Random variable
S_k	One-dimensional sums, equation (30)
T_3	Three-dimensional sum, equation (39)
$U(x)$	Unit step function, equation (5)
\mathbf{x}_n	n -th random variable, equation (1)
z_m	m -th value of parameter, equation (1)
\mathbf{z}_m	m -th ordered random variable
λ	Variable in moment-generating domain, equation (8)
boldface	Random variable

JOINT PROBABILITY DENSITY FUNCTION OF SELECTED ORDER STATISTICS AND THE SUM OF THE REMAINDER AS APPLIED TO ARBITRARY INDEPENDENT RANDOM VARIABLES

INTRODUCTION

Detection and location of weak signals in random noise is frequently accomplished by the ordering of the random variables (RVs) in a measured dataset, followed by an investigation of the locations and statistics of several of the largest RVs under consideration. Also of interest are the remaining smaller RVs in the dataset, which can be used to estimate the background noise level and to form a basis for normalization, thereby realizing a constant false alarm processor.

In this study, the original dataset $\{\mathbf{x}_n\}$ is composed of N independent RVs with arbitrary probability density functions (PDFs) $\{p_n(x)\}$. This dataset is ordered into another dataset of dependent RVs, each with a different PDF. From this latter set, the $M-1$ largest RVs are selected. The sum of the remaining $N+1-M$ RVs is then computed, giving a total of M dependent RVs. The joint M -th order PDF of these M dependent RVs is one of the quantities of interest.

For convenience, the following notation is used. The largest RV in set $\{\mathbf{x}_n\}$ is denoted by \mathbf{z}_1 , the second-largest by \mathbf{z}_2 , ..., the $M-1$ largest by \mathbf{z}_{M-1} , and the sum of the remaining RVs by \mathbf{z}_M . Thus, the first $M-1$ RVs satisfy the restrictions that $\mathbf{z}_1 \geq \mathbf{z}_2 \geq \dots \geq \mathbf{z}_{M-1}$, while the last RV must obviously satisfy the restriction that $\mathbf{z}_M \leq (N+1-M) \mathbf{z}_{M-1}$.

To solve for the statistics of RVs $\{\mathbf{z}_m\}$, general M , a series of simpler problems will be solved starting with $M=2$, that is, the largest RV and the sum of the remainder. By the time M increases to 5, the general pattern will be obvious and may be extended to a larger M of interest. The end result is a single one-dimensional contour integral for the joint PDF of $\{\mathbf{z}_m\}$, which can easily and accurately be numerically evaluated by moving the contour of integration to approximately pass through the real saddlepoint of the integrand.

For later use, it is convenient to define the auxiliary function

$$c_n(z, \lambda) = \int_{-\infty}^z dx \exp(\lambda x) p_n(x) \quad \text{for } n = 1 : N,$$

which is a mixture of a cumulative distribution function (CDF) and a moment-generating function (MGF). That is, $c_n(z, 0)$ is the CDF corresponding to PDF $p_n(x)$, while $c_n(+\infty, \lambda)$ is the corresponding MGF. Variable z is real, while λ can be complex. Several useful examples of the $c_n(z, \lambda)$ function are listed in appendix A of reference 1.

DERIVATIONS OF JOINT PROBABILITY DENSITY FUNCTIONS

The method for deriving joint PDFs is based very heavily on reference 1, pages 5 through 15. The notation introduced there will also be used here. The reader is advised to review that material before proceeding. The major difference here is that each RV now has different statistics. Specifically, the PDF of RV \mathbf{x}_n is $p_n(x)$ for $n = 1:N$, instead of a common PDF $p(x)$ used earlier. The CDF of RV \mathbf{x}_n is $c_n(x)$, while its exceedance distribution function (EDF) is $e_n(x)$. The RVs $\{\mathbf{x}_n\}$, $n = 1:N$, are independent of each other.

LARGEST RANDOM VARIABLE

The probability that RV \mathbf{x}_m is the largest RV and that it is in the interval $(z_1, z_1 + dz_1)$ is

$$dz_1 p_m(z_1) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1). \quad (1)$$

Thus, the probability that RV \mathbf{x}_m is the largest RV is

$$Q(m) \equiv \int dz_1 p_m(z_1) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1) = \int dz \frac{p_m(z)}{c_m(z)} P(z, 0) \quad \text{for } m = 1:N, \quad (2)$$

where the product of auxiliary functions is

$$P(z, \lambda) \equiv \prod_{n=1}^N c_n(z, \lambda). \quad (3)$$

The sum of all the $\{Q(m)\}$ probabilities is unity:

$$\sum_{m=1}^N Q(m) = \int dz \sum_{m=1}^N p_m(z) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z) = \int dz \frac{d}{dz} \prod_{n=1}^N c_n(z) = \prod_{n=1}^N c_n(z) \Big|_{-\infty}^{\infty} = 1. \quad (4)$$

As an example, if all the RVs have exponential PDFs, namely,

$$p_n(x) = a_n \exp(-a_n x) U(x), \quad c_n(x) = [1 - \exp(-a_n x)] U(x) \quad \text{for } n = 1:N, \quad (5)$$

then it follows, from equation (2), that

$$Q(m) = \int_0^{\infty} dz a_m \exp(-a_m z) \prod_{\substack{n=1 \\ n \neq m}}^N \{1 - \exp(-a_n z)\} \quad \text{for } m = 1:N. \quad (6)$$

Although these integrals can be evaluated in closed form, they are probably most efficiently accomplished by numerical integration, especially for large N , once $\{a_n\}$ are specified numerically. For other than exponential PDFs, the integrals in equation (2) will have to be done numerically.

Given that RV $\mathbf{z}_1 = \mathbf{x}_m$ is the largest RV, the conditional PDF of RV \mathbf{x}_n , $n \neq m$, at argument x , when \mathbf{z}_1 has value z_1 , is

$$\frac{p_n(x)}{c_n(z_1)} U(z_1 - x). \quad (7)$$

The corresponding conditional MGF is

$$\int_{-\infty}^{z_1} dx \frac{p_n(x)}{c_n(z_1)} \exp(\lambda x) = \frac{c_n(z_1, \lambda)}{c_n(z_1)}, \quad n \neq m. \quad (8)$$

Therefore, the conditional PDF of the sum \mathbf{z}_2 of the remaining RVs (other than \mathbf{x}_m), at argument z_2 , given that $\mathbf{z}_1 = z_1$, is available from a Bromwich contour integral as

$$\frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \prod_{\substack{n=1 \\ n \neq m}}^N \frac{c_n(z_1, \lambda)}{c_n(z_1)}. \quad (9)$$

Finally, the product of equations (1) and (9) and dz_2 is

$$dz_1 dz_2 p_m(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1, \lambda), \quad (10)$$

which is the probability that $\mathbf{z}_1 (= \mathbf{x}_m)$ is the largest RV, that it lies in the interval $(z_1, z_1 + dz_1)$, and that the sum \mathbf{z}_2 lies in the interval $(z_2, z_2 + dz_2)$. That is,

$$q_2(m, z_1, z_2) \equiv p_m(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1, \lambda) \quad \text{for } m = 1 : N \quad (11)$$

is the combined *probability* (that \mathbf{x}_m is the largest RV) and *joint PDF* of $\mathbf{z}_1 (= \mathbf{x}_m)$ and \mathbf{z}_2 (which is the sum of the remaining RVs).

Let $\lambda = iy$ in equation (11) and integrate on z_2 to get

$$\int dz_2 q_2(m, z_1, z_2) = p_m(z_1) \int dy \delta(y) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1, iy) = p_m(z_1) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1). \quad (12)$$

Then, by use of equation (2), an additional integral yields

$$\iint dz_1 dz_2 q_2(m, z_1, z_2) = Q(m) \text{ for } m = 1:N. \quad (13)$$

That is, function $q_2(m, z_1, z_2)$ in equation (11) is not a true PDF because its area is less than 1. However, the conditional PDF of $\mathbf{z}_1, \mathbf{z}_2$, given that \mathbf{x}_m is the largest RV, is

$$q_2(z_1, z_2 | m) = \frac{1}{Q(m)} p_m(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1, \lambda), \quad (14)$$

which is a true PDF for all $m = 1:N$.

Alternatively, the sum of equation (11) over all m ,

$$q_2(z_1, z_2) \equiv \sum_{m=1}^N q_2(m, z_1, z_2) = \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \sum_{m=1}^N p_m(z_1) \prod_{\substack{n=1 \\ n \neq m}}^N c_n(z_1, \lambda), \quad (15)$$

is a true PDF, namely, the unconditional joint PDF of the largest RV \mathbf{z}_1 and the sum of the remaining RVs \mathbf{z}_2 . Using equation (3), this quantity can be expressed as

$$q_2(z_1, z_2) = \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) P(z_1, \lambda) \sum_{n=1}^N \frac{p_n(z_1)}{c_n(z_1, \lambda)}, \quad (16)$$

which avoids the nested two-dimensional operation in equation (15) in favor of a one-dimensional product and a one-dimensional sum, both of which depend on the variable of integration λ . This joint PDF is zero for $z_2 > (N-1)z_1$ because $\text{RV } \mathbf{z}_2 < (N-1)\mathbf{z}_1$ is always true. An alternative argument is given in equation (16) of reference 1; it uses the fact that $c_n(z, \lambda)$ is analytic in λ for $\text{Re}(\lambda) > \lambda_n$, a problem-dependent critical value.

As an alternative check, if all the RVs $\{\mathbf{x}_n\}$ are identically distributed, equation (15) reduces to equation (15) of reference 1.

TWO LARGEST RANDOM VARIABLES

The probability that RV \mathbf{x}_m is the largest RV, that \mathbf{x}_k is the second-largest RV, that \mathbf{x}_m is in interval $(z_1, z_1 + dz_1)$, and that \mathbf{x}_k is in interval $(z_2, z_2 + dz_2)$, $z_1 > z_2$, is

$$dz_1 dz_2 p_m(z_1) p_k(z_2) \prod_{\substack{n=1 \\ n \neq m, k}}^N c_n(z_2) U(z_1 - z_2), \quad m \neq k. \quad (17)$$

Then, the probability that \mathbf{x}_m is the largest RV and that \mathbf{x}_k is the second-largest RV is best obtained by integrating equation (17) over z_1 first:

$$\begin{aligned}
Q(m,k) &\equiv \int dz_2 p_k(z_2) \prod_{\substack{n=1 \\ n \neq m,k}}^N c_n(z_2) \int_{z_2}^{\infty} dz_1 p_m(z_1) \\
&= \int dz_2 e_m(z_2) p_k(z_2) \prod_{\substack{n=1 \\ n \neq m,k}}^N c_n(z_2) \\
&= \int dz P(z,0) \frac{e_m(z) p_k(z)}{c_m(z) c_k(z)} \quad \text{for } m \neq k, \quad m, k = 1 : N.
\end{aligned} \tag{18}$$

A single integral suffices to determine this probability. The sum of all the $\{Q(m,k)\}$ probabilities for $m \neq k$ is unity.

For the example of exponential RVs in equation (5), the integral in equation (18) takes the form

$$Q(m,k) = \int_0^{\infty} dz a_k \exp[-(a_m + a_k)z] \prod_{\substack{n=1 \\ n \neq m,k}}^N \{1 - \exp(-a_n z)\} \quad \text{for } m \neq k. \tag{19}$$

Again, although possible analytically, numerical integration is the most practical method.

Given that RV $\mathbf{z}_1 = \mathbf{x}_m$ is the largest RV and that $\mathbf{z}_2 = \mathbf{x}_k$ is the second-largest RV, the conditional PDF of \mathbf{x}_n , $n \neq m, k$, at argument x , when \mathbf{z}_1 has value z_1 and \mathbf{z}_2 has value z_2 , with $z_1 > z_2$, is

$$\frac{p_n(x)}{c_n(z_2)} U(z_2 - x). \tag{20}$$

The corresponding conditional MGF is

$$\int_{-\infty}^{z_2} dx \frac{p_n(x)}{c_n(z_2)} \exp(\lambda x) = \frac{c_n(z_2, \lambda)}{c_n(z_2)}, \quad n \neq m, k. \tag{21}$$

The conditional PDF of the sum \mathbf{z}_3 of the remaining RVs (other than $\mathbf{x}_m, \mathbf{x}_k$) at argument z_3 , given that $\mathbf{z}_1 = z_1$ and $\mathbf{z}_2 = z_2$, is

$$\frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) \prod_{\substack{n=1 \\ n \neq m,k}}^N \frac{c_n(z_2, \lambda)}{c_n(z_2)}. \tag{22}$$

The product of equations (17) and (22) and dz_3 is

$$dz_1 dz_2 dz_3 p_m(z_1) p_k(z_2) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) \prod_{\substack{n=1 \\ n \neq m, k}}^N c_n(z_2, \lambda) U(z_1 - z_2), \quad m \neq k, \quad (23)$$

which is the probability that \mathbf{x}_m is the largest RV, that \mathbf{x}_k is the second-largest RV, that $\mathbf{z}_1 = \mathbf{x}_m$ lies in $(z_1, z_1 + dz_1)$, that $\mathbf{z}_2 = \mathbf{x}_k$ lies in $(z_2, z_2 + dz_2)$, and that sum \mathbf{z}_3 lies in $(z_3, z_3 + dz_3)$. That is, for $m \neq k$,

$$q_3(m, k, z_1, z_2, z_3) \equiv U(z_1 - z_2) p_m(z_1) p_k(z_2) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) \prod_{\substack{n=1 \\ n \neq m, k}}^N c_n(z_2, \lambda) \quad (24)$$

is the combined probability (that \mathbf{x}_m is the largest RV and that \mathbf{x}_k is the second-largest RV) and joint PDF of $\mathbf{z}_1 (= \mathbf{x}_m)$, $\mathbf{z}_2 (= \mathbf{x}_k)$, and \mathbf{z}_3 (the sum of the remaining RVs).

Let $\lambda = iy$ in equation (24) and integrate on z_3 to get

$$\int dz_3 q_3(m, k, z_1, z_2, z_3) = U(z_1 - z_2) p_m(z_1) p_k(z_2) \prod_{\substack{n=1 \\ n \neq m, k}}^N c_n(z_2). \quad (25)$$

Then, by reference to equations (17) and (18), the remaining two integrals yield

$$\iiint dz_1 dz_2 dz_3 q_3(m, k, z_1, z_2, z_3) = Q(m, k), \quad m \neq k. \quad (26)$$

That is, function $q_3(m, k, z_1, z_2, z_3)$ in equation (24) is not a true PDF because its area is less than 1. However, the conditional PDF of $\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3$, *given* that \mathbf{z}_m is the largest RV and that \mathbf{z}_k is the second-largest RV, is

$$\begin{aligned} q_3(z_1, z_2, z_3 | m, k) &= \frac{1}{Q(m, k)} U(z_1 - z_2) p_m(z_1) p_k(z_2) \\ &\times \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) \prod_{\substack{n=1 \\ n \neq m, k}}^N c_n(z_2, \lambda), \end{aligned} \quad (27)$$

which is a true PDF for all $m \neq k$, $m, k = 1:N$.

Alternatively, the sum of equation (24) over all $m \neq k$,

$$\begin{aligned}
 q_3(z_1, z_2, z_3) &\equiv \sum_{\substack{m,k=1 \\ m \neq k}}^N q_3(m, k, z_1, z_2, z_3) \\
 &= U(z_1 - z_2) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) \sum_{m=1}^N p_m(z_1) \sum_{\substack{k=1 \\ k \neq m}}^N p_k(z_2) \prod_{\substack{n=1 \\ n \neq m, k}}^N c_n(z_2, \lambda),
 \end{aligned} \tag{28}$$

is a true PDF, namely, the unconditional PDF of the largest RV z_1 , the second-largest RV z_2 , and the sum of the remaining RVs z_3 . By using equation (3), and adding and subtracting the $k = m$ term in the inner sum, this quantity can be expressed as

$$q_3(z_1, z_2, z_3) = U(z_1 - z_2) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) P(z_2, \lambda) (S_1 S_2 - S_3), \tag{29}$$

where one-dimensional sums

$$S_1 = \sum_{n=1}^N \frac{p_n(z_1)}{c_n(z_2, \lambda)}, \quad S_2 = \sum_{n=1}^N \frac{p_n(z_2)}{c_n(z_2, \lambda)}, \quad S_3 = \sum_{n=1}^N \frac{p_n(z_1) p_n(z_2)}{c_n(z_2, \lambda)^2}. \tag{30}$$

Equation (29) is much more advantageous computationally than equation (28), which requires a nested three-dimensional sum and product. Joint PDF $q_3(z_1, z_2, z_3)$ is zero for $z_3 > (N-2) z_2$ because RV $z_3 < (N-2) z_2$ is always true.

As a check, if all the RVs $\{x_n\}$ are identically distributed, the result in equation (28) reduces to

$$N(N-1) p(z_1) p(z_2) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_3) c(z_2, \lambda)^{N-2} U(z_1 - z_2), \tag{31}$$

which is equation (35) of reference 1. Also, by letting $\lambda = iy$ in equation (31) and integrating on z_3 , there follows

$$N(N-1) p(z_1) p(z_2) c(z_2)^{N-2} U(z_1 - z_2). \tag{32}$$

An additional integral on z_2 yields

$$N(N-1) p(z_1) \int_{-\infty}^{z_1} dz_2 p(z_2) c(z_2)^{N-2} = N p(z_1) c(z_1)^{N-1} = \frac{d}{dz_1} c(z_1)^N, \tag{33}$$

where the step function $U(z_1 - z_2)$ takes effect. Finally, integrating on z_1 yields 1.

THREE LARGEST RANDOM VARIABLES

The probability that RV \mathbf{x}_m is the largest RV and lies in interval $(z_1, z_1 + dz_1)$, that \mathbf{x}_k is the second-largest RV and lies in $(z_2, z_2 + dz_2)$, and that \mathbf{x}_j is the third-largest RV and lies in $(z_3, z_3 + dz_3)$, with $z_1 > z_2 > z_3$, is

$$dz_1 dz_2 dz_3 p_m(z_1) p_k(z_2) p_j(z_3) \prod_{\substack{n=1 \\ n \neq m, k, j}}^N c_n(z_3) U(z_1 - z_2) U(z_2 - z_3) \quad (34)$$

for m, k, j all unequal. Then, the probability that \mathbf{x}_m is the largest RV, that \mathbf{x}_k is the second-largest RV, and that \mathbf{x}_j is the third-largest RV is obtained by integrating equation (34) first on z_1 to obtain

$$Q(m, k, j) = \int_{-\infty}^{\infty} dz_3 p_j(z_3) \prod_{\substack{n=1 \\ n \neq m, k, j}}^N \{c_n(z_3)\} \int_{z_3}^{\infty} dz_2 p_k(z_2) e_m(z_2). \quad (35)$$

At this point, in general, the remaining double integral cannot be reduced any further, although the sum of all the $\{Q(m, k, j)\}$ over all unequal m, k, j must be unity. However, for the example of exponential RVs in equation (5), the z_2 integral can be carried out to yield

$$Q(m, k, j) = \frac{a_k a_j}{a_m + a_k} \int_0^{\infty} dz \exp[-(a_m + a_k + a_j) z] \prod_{\substack{n=1 \\ n \neq m, k, j}}^N \{1 - \exp(-a_n z)\}. \quad (36)$$

For given numerical values of $\{a_n\}$, this single integral can be easily evaluated for any m, k, j of interest.

For general statistics of RVs $\{\mathbf{x}_n\}$, and by a similar procedure to that presented in equations (20) through (24), the combined probability and joint PDF of $\mathbf{z}_1 (= \mathbf{x}_m)$, $\mathbf{z}_2 (= \mathbf{x}_k)$, $\mathbf{z}_3 (= \mathbf{x}_j)$, and \mathbf{z}_4 (the sum of the remaining RVs) is

$$q_4(m, k, j, z_1, z_2, z_3, z_4) = U(z_1 - z_2) U(z_2 - z_3) p_m(z_1) p_k(z_2) p_j(z_3) \\ \times \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_4) \prod_{\substack{n=1 \\ n \neq m, k, j}}^N c_n(z_3, \lambda) \quad \text{for } m, k, j \text{ all unequal.} \quad (37)$$

Function q_4 is zero for $z_4 > (N - 3) z_3$ because RV $\mathbf{z}_4 < (N - 3) \mathbf{z}_3$ is always true. Equation (37) can be evaluated numerically with moderate computational effort.

The sum of equation (37) over all unequal m, k, j is the unconditional joint PDF of RVs $\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3$, and \mathbf{z}_4 and can be expressed as

$$q_4(z_1, z_2, z_3, z_4) = U(z_1 - z_2) U(z_2 - z_3) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_4) P(z_3, \lambda) T_3, \quad (38)$$

where

$$T_3 = T(z_1, z_2, z_3, \lambda) = \sum_{m=1}^N \frac{p_m(z_1)}{c_m(z_3, \lambda)} \sum_{\substack{k=1 \\ k \neq m}}^N \frac{p_k(z_2)}{c_k(z_3, \lambda)} \sum_{\substack{j=1 \\ j \neq m, k}}^N \frac{p_j(z_3)}{c_j(z_3, \lambda)}. \quad (39)$$

By adding and subtracting the missing terms in each sum, starting with the innermost sum and expanding out the resulting expressions, the following form is obtained:

$$T_3 = S_1 S_2 S_3 - S_1 S_6 - S_2 S_5 - S_3 S_4 + 2S_7, \quad (40)$$

where

$$S_1 = \sum_n a_n, \quad S_2 = \sum_n b_n, \quad S_3 = \sum_n c_n, \quad S_4 = \sum_n a_n b_n, \quad (41)$$

$$S_5 = \sum_n a_n c_n, \quad S_6 = \sum_n b_n c_n, \quad S_7 = \sum_n a_n b_n c_n,$$

and

$$a_n = \frac{p_n(z_1)}{c_n(z_3, \lambda)}, \quad b_n = \frac{p_n(z_2)}{c_n(z_3, \lambda)}, \quad c_n = \frac{p_n(z_3)}{c_n(z_3, \lambda)} \quad \text{for } n = 1 : N. \quad (42)$$

Whereas direct evaluation of equation (39) would require a triple-nested sum, requiring a number of operations of the order of N^3 , equation (40) requires only the seven one-dimensional sums in equation (41), each of size N .

The PDF $q_4(z_1, z_2, z_3, z_4)$ in equation (38) is zero for $z_4 > (N-3)z_3$ because RV $\mathbf{z}_4 < (N-3)\mathbf{z}_3$ is always true.

FOUR LARGEST RANDOM VARIABLES

The probability that RV \mathbf{x}_m is the largest RV and lies in interval $(z_1, z_1 + dz_1)$, that \mathbf{x}_k is the second-largest RV and lies in $(z_2, z_2 + dz_2)$, that \mathbf{x}_j is the third-largest RV and lies in $(z_3, z_3 + dz_3)$, and that \mathbf{x}_i is the fourth-largest RV and lies in $(z_4, z_4 + dz_4)$, where $z_1 > z_2 > z_3 > z_4$, is

$$dz_1 dz_2 dz_3 dz_4 p_m(z_1) p_k(z_2) p_j(z_3) p_i(z_4) \prod_{\substack{n=1 \\ n \neq m, k, j, i}}^N c_n(z_4) U(z_1 - z_2) U(z_2 - z_3) U(z_3 - z_4) \quad (43)$$

for m, k, j, i all unequal. Then, the probability that \mathbf{x}_m is the largest RV, that \mathbf{x}_k is the second-largest RV, that \mathbf{x}_j is the third-largest RV, and that \mathbf{x}_i is the fourth-largest RV is obtained by integrating equation (43) first on z_1 to obtain

$$Q(m, k, j, i) = \int_{-\infty}^{\infty} dz_4 p_i(z_4) \prod_{\substack{n=1 \\ n \neq m, k, j, i}}^N \{c_n(z_4)\} \int_{z_4}^{\infty} dz_3 p_j(z_3) \int_{z_3}^{\infty} dz_2 p_k(z_2) e_m(z_2). \quad (44)$$

At this point, in general, the remaining triple integral cannot be reduced any further, although the sum of all the $\{Q(m, k, j, i)\}$ over all unequal m, k, j, i must be unity. However, for the example of exponential RVs in equation (5), the z_2 and z_3 integrals can be carried out to yield

$$Q(m, k, j, i) = \frac{a_k a_j a_i}{(a_m + a_k)(a_m + a_k + a_j)} \times \int_0^{\infty} dz \exp[-(a_m + a_k + a_j + a_i) z] \prod_{\substack{n=1 \\ n \neq m, k, j, i}}^N \{1 - \exp(-a_n z)\}. \quad (45)$$

For given numerical values of $\{a_n\}$, this single integral can be easily evaluated for any m, k, j, i of interest.

For general statistics of RVs $\{\mathbf{x}_n\}$, and using a procedure similar to that presented in equations (20) through (24), the combined probability and joint PDF of $\mathbf{z}_1 (= \mathbf{x}_m)$, $\mathbf{z}_2 (= \mathbf{x}_k)$, $\mathbf{z}_3 (= \mathbf{x}_j)$, $\mathbf{z}_4 (= \mathbf{x}_i)$, and \mathbf{z}_5 (the sum of the remaining RVs) is

$$q_5(m, k, j, i, z_1, z_2, z_3, z_4, z_5) = U(z_1 - z_2) U(z_2 - z_3) U(z_3 - z_4) p_m(z_1) p_k(z_2) \times p_j(z_3) p_i(z_4) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_5) \prod_{\substack{n=1 \\ n \neq m, k, j, i}}^N c_n(z_4, \lambda) \quad (46)$$

for m, k, j, i all unequal.

Function q_5 is zero for $z_5 > (N-4)z_4$ because RV $\mathbf{z}_5 < (N-4)\mathbf{z}_4$ is always true. Equation (46) can be evaluated numerically with moderate computational effort.

The sum of equation (46) over all unequal m, k, j, i is the unconditional joint PDF of RVs $\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3, \mathbf{z}_4$, and \mathbf{z}_5 and is given by

$$q_5(z_1, \dots, z_5) = U(z_1 - z_2) U(z_2 - z_3) U(z_3 - z_4) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_5) P(z_4, \lambda) T_4, \quad (47)$$

where, after expansion and manipulation, the fourth-order sum for T_4 can be expressed as

$$\begin{aligned} T_4 = T(z_1, z_2, z_3, z_4, \lambda) = & S_1 S_2 S_3 S_4 + 2(S_1 S_{14} + S_2 S_{13} + S_3 S_{12} + S_4 S_{11}) \\ & + S_5 S_{10} + S_6 S_9 + S_7 S_8 - S_1 (S_2 S_{10} + S_3 S_9 + S_4 S_8) - S_2 (S_3 S_7 + S_4 S_6) \\ & - S_3 S_4 S_5 - 6S_{15}, \end{aligned} \quad (48)$$

and the 15 sums are given by

$$\begin{aligned} S_1 &= \sum_n a_n, \quad S_2 = \sum_n b_n, \quad S_3 = \sum_n c_n, \quad S_4 = \sum_n d_n, \\ S_5 &= \sum_n a_n b_n, \quad S_6 = \sum_n a_n c_n, \quad S_7 = \sum_n a_n d_n, \quad S_8 = \sum_n b_n c_n, \\ S_9 &= \sum_n b_n d_n, \quad S_{10} = \sum_n c_n d_n, \quad S_{11} = \sum_n a_n b_n c_n, \quad S_{12} = \sum_n a_n b_n d_n, \\ S_{13} &= \sum_n a_n c_n d_n, \quad S_{14} = \sum_n b_n c_n d_n, \quad S_{15} = \sum_n a_n b_n c_n d_n, \end{aligned} \quad (49)$$

with

$$a_n = \frac{p_n(z_1)}{c_n(z_4, \lambda)}, \quad b_n = \frac{p_n(z_2)}{c_n(z_4, \lambda)}, \quad c_n = \frac{p_n(z_3)}{c_n(z_4, \lambda)}, \quad d_n = \frac{p_n(z_4)}{c_n(z_4, \lambda)} \quad \text{for } n = 1:N. \quad (50)$$

Again, the simplification afforded by form (48) is a considerable improvement on the initial quadruple nested sum encountered for T_4 in equation (47). All 15 sums in equation (49) are of size N . An additional shortcut is available by defining $e_n = a_n b_n$ and $f_n = c_n d_n$ for $n = 1:N$ in all the sums in equation (49).

The PDF $q_5(z_1, z_2, z_3, z_4, z_5)$ in equation (47) is zero for $z_5 > (N-4)z_4$ because RV $\mathbf{z}_5 < (N-4)\mathbf{z}_4$ is always true.

M-1 LARGEST RANDOM VARIABLES

Let $\mathbf{x}_{m_1} = \mathbf{z}_1$ be the largest RV, $\mathbf{x}_{m_2} = \mathbf{z}_2$ be the second-largest RV, and $\mathbf{x}_{m_{M-1}} = \mathbf{z}_{M-1}$ be the $(M-1)$ -th-largest RV. Also, let \mathbf{z}_M be the sum of the remaining RVs. Then, $\mathbf{z}_M < (N+1-M) \mathbf{z}_{M-1}$ is always true.

Observation of equation (46) reveals that the combined probability and joint PDF of the M RVs $\{\mathbf{z}_m\}$, $m = 1:M$, is given by

$$q_M(m_1, \dots, m_{M-1}; z_1, \dots, z_M) = \prod_{m=1}^{M-2} \{U(z_m - z_{m+1})\} \prod_{j=1}^{M-1} \{p_{m_j}(z_j)\} \times \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_M) P(z_{M-1}, \lambda) / \prod_{j=1}^{M-1} \{c_{m_j}(z_{M-1}, \lambda)\}. \quad (51)$$

For this quantity to be nonzero, it is required that

$$z_1 > z_2 > \dots > z_{M-1} \text{ and } z_M < (N+1-M) z_{M-1}. \quad (52)$$

RECURSION FOR m -TH ORDER SUM

As $M-1$, the number of the largest RVs of interest, increases, the initial form for the nested sum $T_m \equiv T(z_1, \dots, z_m, \lambda)$, $m = M-1$, becomes impractical computationally. Also, the expansion and simplification procedure leading to compact equations (40) and (48) for $m = 3$ and $m = 4$, respectively, becomes very tedious and unwieldy. A method around these limitations is to develop a recursion procedure for getting T_m directly from T_{m-1} .

The essentials of this derivation begin with the definition

$$A_2(a, b) \equiv \sum_{n=1}^N a_n \sum_{\substack{m=1 \\ m \neq n}}^N b_m = \sum_{n=1}^N a_n \sum_{m=1}^N b_m - \sum_{n=1}^N a_n b_n = \text{sum}(a) * \text{sum}(b) - \text{sum}(a * b). \quad (53)$$

Suppose a program is written to perform this task on sequences $\{a_n\}$ and $\{b_n\}$. Now, consider the third-order nested sum

$$A_3(a, b, c) \equiv \sum_{n=1}^N a_n \sum_{\substack{m=1 \\ m \neq n}}^N b_m \sum_{\substack{k=1 \\ k \neq m, n}}^N c_k. \quad (54)$$

Develop the inner sum according to

$$\begin{aligned}
 A_3(a, b, c) &= \sum_{n=1}^N a_n \sum_{\substack{m=1 \\ m \neq n}}^N b_m \left(\sum_{k=1}^N c_k - c_n - c_m \right) \\
 &= \text{sum}(c) * A_2(a, b) - A_2(a, *c, b) - A_2(a, b, *c),
 \end{aligned} \tag{55}$$

where the notation introduced in equation (53) has been used. Thus, A_3 can be evaluated by three calls to function A_2 . It follows in a similar fashion that

$$\begin{aligned}
 A_4(a, b, c, d) &\equiv \sum_{n=1}^N a_n \sum_{\substack{m=1 \\ m \neq n}}^N b_m \sum_{\substack{k=1 \\ k \neq m, n}}^N c_k \sum_{\substack{j=1 \\ j \neq k, m, n}}^N d_j \\
 &= \text{sum}(d) * A_3(a, b, c) - A_3(a, *d, b, c) - A_3(a, b, *d, c) - A_3(a, b, c, *d).
 \end{aligned} \tag{56}$$

That is, in general, A_m can be evaluated by m calls to A_{m-1} , with appropriate combinations of arguments.

It should be observed that the use of the recursive approach is not as economical as having explicit expressions for the high-order sums of interest. For example, the use of equation (55) employs 10 summations, whereas the direct use of equation (40) requires that only seven sums be evaluated. (Sums S_1, S_2 , and S_7 are evaluated twice in equation (55).) This loss of economy is present at every level and gets worse as m increases. Also, the amount of computational effort increases noticeably with m . In fact, the number of summations that must be evaluated at level m is $2^m - 1$; observe the results in equations (30), (41), and (49), for example. This additional effort serves to effectively limit the level to which the procedure can be carried out; that is, evaluation of the joint PDF of the m largest RVs and the sum of the remainder is not computationally feasible for very large m . Additional effort on extending results like equations (48) through (50) would probably be very worthwhile, at least for $m = 5$ or 6. For example, at $m = 5$, the number of different sums to be evaluated is $2^5 - 1 = 31$, and the number of different types of terms in T_5 is 52. Namely,

$$T_5 = 24 Q(5) - 6 Q(4,1) - 2 Q(3,2) + 2 Q(3,1,1) + Q(2,2,1) - Q(2,1,1,1) + Q(1,1,1,1,1), \tag{57}$$

where

$$Q(5) = \sum a b c d e, \text{ 1 term,}$$

$$Q(4,1) = \sum a b c d \sum e + \dots, \text{ 5 terms,}$$

$$Q(3,2) = \sum a b c \sum d e + \dots, \text{ 10 terms,}$$

$$Q(3,1,1) = \sum a b c \sum d \sum e + \dots, \text{ 10 terms,} \quad (58)$$

$$Q(2,2,1) = \sum a b \sum c d \sum e + \dots, \text{ 15 terms,}$$

$$Q(2,1,1,1) = \sum a b \sum c \sum d \sum e + \dots, \text{ 10 terms,}$$

$$Q(1,1,1,1,1) = \sum a \sum b \sum c \sum d \sum e, \text{ 1 term.}$$

SECOND-LARGEST RANDOM VARIABLE

The probability that \mathbf{x}_k is the largest RV, that \mathbf{x}_j is the second-largest RV, and that $\mathbf{x}_j \in (z_1, z_1 + dz_1)$ is

$$dz_1 p_j(z_1) e_k(z_1) \prod_{\substack{n=1 \\ n \neq j, k}}^N c_n(z_1), \quad k \neq j. \quad (59)$$

Given that $\mathbf{x}_j = z_1$, the conditional PDF of \mathbf{x}_k is $\frac{p_k(x)}{e_k(z_1)} U(x - z_1)$, while that for \mathbf{x}_n , $n \neq j, k$, is

$\frac{p_n(x)}{c_n(z_1)} U(z_1 - x)$. The corresponding conditional MGF of \mathbf{x}_k is

$$\int_{z_1}^{\infty} dx \frac{p_k(x)}{e_k(z_1)} \exp(\lambda x) = \frac{e_k(z_1, \lambda)}{e_k(z_1)} \quad (60)$$

and that of \mathbf{x}_n is

$$\int_{-\infty}^{z_1} dx \frac{p_n(x)}{c_n(z_1)} \exp(\lambda x) = \frac{c_n(z_1, \lambda)}{c_n(z_1)}, \quad n \neq j, k. \quad (61)$$

The conditional MFG of $\mathbf{z}_2 = \mathbf{x}_k + \sum_{\substack{n=1 \\ n \neq j, k}}^N \mathbf{x}_n$ is

$$\frac{e_k(z_1, \lambda)}{e_k(z_1)} \prod_{\substack{n=1 \\ n \neq j, k}}^N \frac{c_n(z_1, \lambda)}{c_n(z_1)}. \quad (62)$$

The conditional PDF of \mathbf{z}_2 at argument z_2 , given $\mathbf{z}_1 = \mathbf{x}_j$ has value z_1 , is

$$\frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \frac{e_k(z_1, \lambda)}{e_k(z_1)} \prod_{\substack{n=1 \\ n \neq j, k}}^N \frac{c_n(z_1, \lambda)}{c_n(z_1)}. \quad (63)$$

Finally, the product of equations (59) and (63) and dz_2 is

$$dz_1 dz_2 p_j(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) e_k(z_1, \lambda) \prod_{\substack{n=1 \\ n \neq j, k}}^N c_n(z_1, \lambda), \quad k \neq j, \quad (64)$$

which is the probability that \mathbf{x}_k is the largest RV, that \mathbf{x}_j is the second-largest RV, that $\mathbf{z}_1 = \mathbf{x}_j \in (z_1, z_1 + dz_1)$, and that $\mathbf{z}_2 \in (z_2, z_2 + dz_2)$. That is,

$$q_2(j, k, z_1, z_2) = p_j(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) e_k(z_1, \lambda) \prod_{\substack{n=1 \\ n \neq j, k}}^N c_n(z_1, \lambda), \quad k \neq j, \quad (65)$$

is the combined probability (that \mathbf{x}_k is the largest RV and that \mathbf{x}_j is the second-largest RV) and joint PDF of \mathbf{z}_1 , the second-largest RV, and \mathbf{z}_2 , the sum of the remaining RVs.

The sum of equation (65) over k ,

$$q_2(j, z_1, z_2) = p_j(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \sum_{\substack{k=1 \\ k \neq j}}^N e_k(z_1, \lambda) \prod_{\substack{n=1 \\ n \neq j, k}}^N c_n(z_1, \lambda), \quad (66)$$

is the combined probability (that \mathbf{x}_j is the second-largest RV) and joint PDF of $\mathbf{z}_1 (= \mathbf{x}_j)$ and \mathbf{z}_2 (the sum of the remaining RVs). This quantity can be written as

$$q_2(j, z_1, z_2) = p_j(z_1) \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) \frac{P(z_1, \lambda)}{c_j(z_1, \lambda)} \sum_{\substack{k=1 \\ k \neq j}}^N \frac{e_k(z_1, \lambda)}{c_k(z_1, \lambda)} \quad \text{for } j = 1 : N. \quad (67)$$

Then, the sum of equation (67) over j ,

$$q_2(z_1, z_2) = \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) P(z_1, \lambda) \sum_{j=1}^N \frac{p_j(z_1)}{c_j(z_1, \lambda)} \sum_{\substack{k=1 \\ k \neq j}}^N \frac{e_k(z_1, \lambda)}{c_k(z_1, \lambda)}, \quad (68)$$

is the overall joint PDF of the second-largest RV z_1 and the sum of the remaining RVs z_2 . This last quantity can be expressed alternatively as

$$q_2(z_1, z_2) = \frac{1}{i2\pi} \int_C d\lambda \exp(-\lambda z_2) P(z_1, \lambda) (S_1 S_2 - S_3), \quad (69)$$

where one-dimensional sums

$$S_1 = \sum_{n=1}^N \frac{p_n(z_1)}{c_n(z_1, \lambda)}, \quad S_2 = \sum_{n=1}^N \frac{e_n(z_1, \lambda)}{c_n(z_1, \lambda)}, \quad S_3 = \sum_{n=1}^N \frac{p_n(z_1) e_n(z_1, \lambda)}{c_n(z_1, \lambda)^2}. \quad (70)$$

SUMMARY

The joint statistics of $M-1$ ordered random variables and the sum of the remaining random variables have been derived for several values of low-order M . The original random variables, prior to ordering, are independent and can have arbitrary, different probability density functions. Results for the joint probability density function of the M random variables of interest, as well as a combined probability and joint probability density function, have been derived in the form of a single contour integral in the moment-generating domain. Numerical evaluation of this contour integral is most easily accomplished by approximately locating the real saddlepoint of the integrand and moving the Bromwich contour so as to pass through this point. However, instead of resorting to a saddlepoint approximation, high accuracy in the evaluation of the joint probability density function is achievable by numerical integration along this displaced contour.

A recursive procedure has been developed for evaluating a nested sum that occurs in the evaluation of the joint probability density function. For values of M in the range of 6 to 10, this is a very helpful numerical aid. For much larger values of M , execution time increases very rapidly and becomes a significant limitation.

REFERENCES

1. A. H. Nuttall, "Joint Probability Density Function of Selected Order Statistics and the Sum of the Remaining Random Variables," NUWC-NPT Technical Report 11,345, Naval Undersea Warfare Center Division, Newport, RI, 15 January 2002.

INITIAL DISTRIBUTION LIST

Addressee	No. of Copies
Office of Naval Research (D. Abraham, J. Tague, B. Fitch)	3
University of Connecticut (P. Willett)	1
University of Rhode Island (S. Kay)	1
Defense Technical Information Center	2